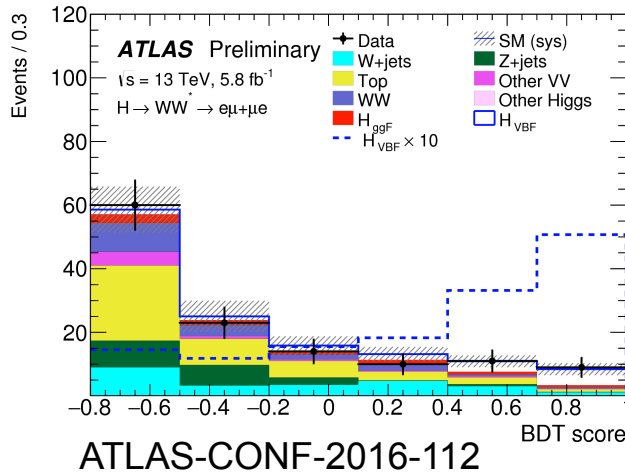# Weakly supervised classifiers
## learning from data and proportions
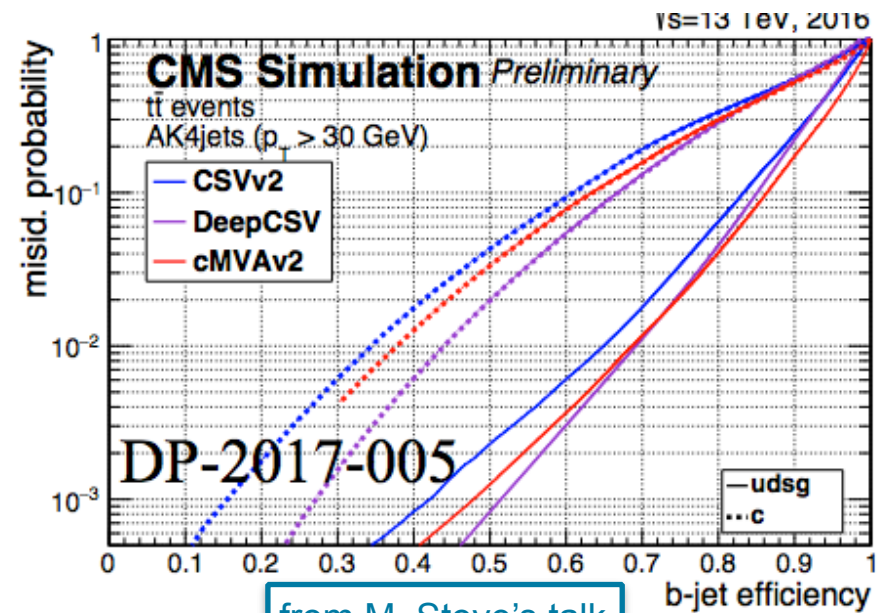
L. Dery (Stanford), B. Nachman (LBNL), F. Rubbo (SLAC), A. Schwartzman (SLAC)

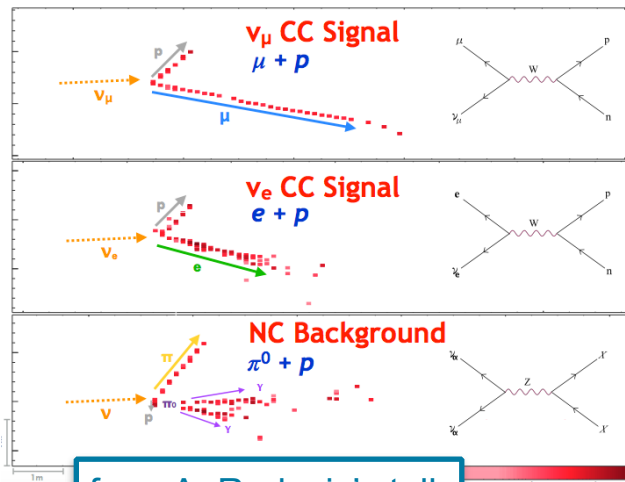# Classification in HEP

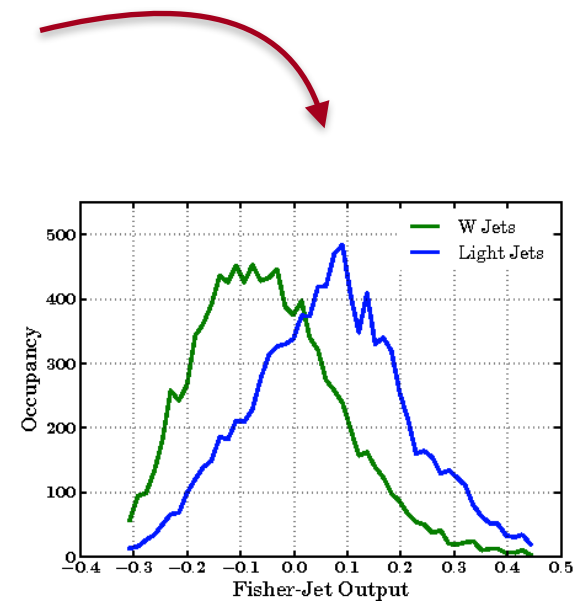Discriminating signal events from backgrounds



ATLAS-CONF-2016-112



from A. Radovic's talk

Classifying reconstructed objects



from M. Stoye's talk

& more…

# Jet classification example

CNN
RNN

W Jets
Light Jets

Occupancy

Fisher-Jet Output

# Jet classification example

1511.05190
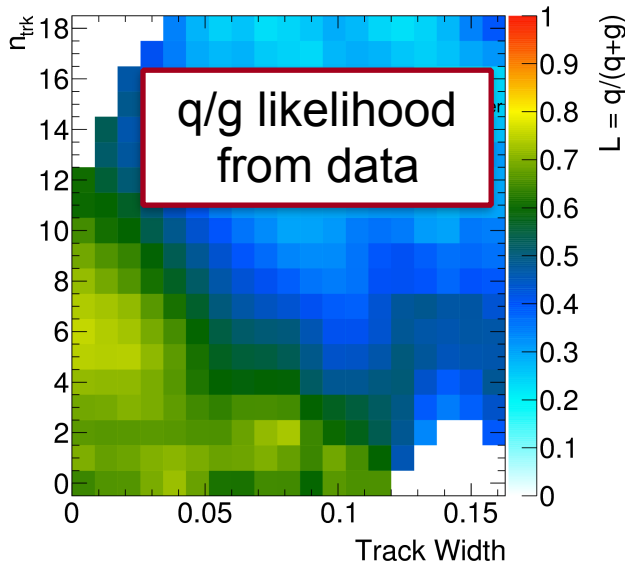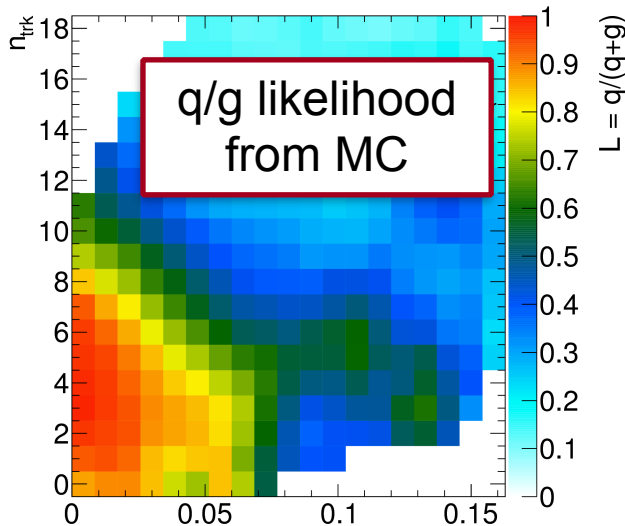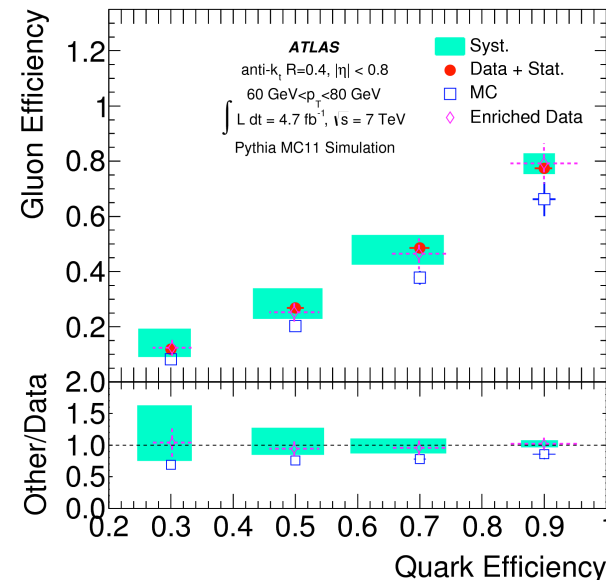


*or other high-dimensional representations (embedding,M-body, etc..)

# Learning from simulation vs learning from data



q/g likelihood from MC



q/g likelihood from data

- Modeling of multi-dimensional soft QCD features (e.g. $n_{track}$, $w_{track}$) is challenging for MC.



1405.6583

- Expect further strain at higher dimensionality (e.g. images with thousands of pixels!)

5

# Training on data

- Classifier is always suboptimal if distribution of training and test samples are different.

- Data is the perfect event "simulation": exactly the same distribution as in the test sample.

- N.B.: doesn't impact uncertainties, only the "central value" of the performance (i.e. how optimal is the discrimination in data)!

- N.B.2: for many applications simulation is very good and its distribution is close to data.

# Learn directly from unlabeled data!



**Weakly supervised classifier trained without using labels**

**Traditional fully supervised classifier**

Legend:
- Weakly supervised NN, AUC=0.93
- Fully supervised NN, AUC=0.93
- Feature 1, auc=0.77
- Feature 2, auc=0.70
- Feature 3, auc=0.78
- Feature 4, auc=0.77
- Feature 5, auc=0.70

# Traditional full supervision

Labeled training set ("simulation")

apples

pears

$$f_{\text{full}} = \text{argmin}_{f':\mathbb{R}^n \to \{0,1\}} \sum_{i=1}^{N} \ell\left(f'(x_i) - t_i\right)$$

instance label:
**0:pear 1:apple**

Classification $f_{\text{full}}\left( \right) = 0.97$

# Weak supervision

unlabeled training data

2/3 apples

1/3 apples

2/5 apples

1/4 apples

$$f_{\text{weak}} = \operatorname{argmin}_{f':\mathbb{R}^n \to [0,1]} \ell\left(\sum_{i=1}^{N} \frac{f'(x_i)}{N} - y\right)$$

average composition for each barrel

Classification $f_{\text{weak}}$ ( 🍎 ) = 0.97

# Weak supervision - analytically

unlabeled data sample A

$y_A = 0.1$

i$^{th}$ bin   X

unlabeled data sample B

$y_B = 0.3$

i$^{th}$ bin   X

$$h_{A,i} = y_A h_{1,i} + (1 - y_A) h_{0,i}$$

$$h_{B,i} = y_B h_{1,i} + (1 - y_B) h_{0,i}$$

- Given two independent unlabeled data samples, and the corresponding proportion of signal, we can extract the signal and background distributions.

# Weak supervision - analytically

unlabeled data sample A

**signal**

**background**

unlabeled data sample B

$y_A = 0.1$

$y_B = 0.3$

X
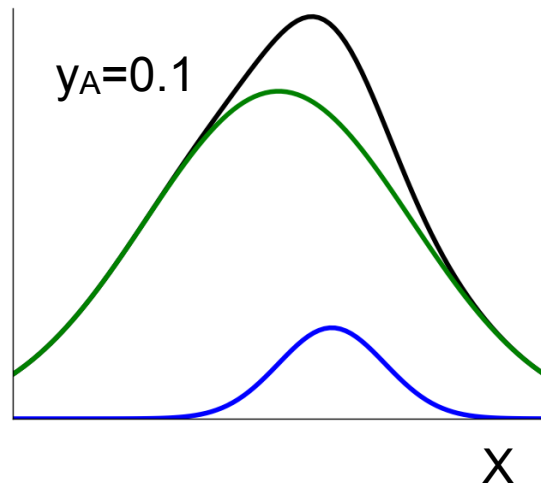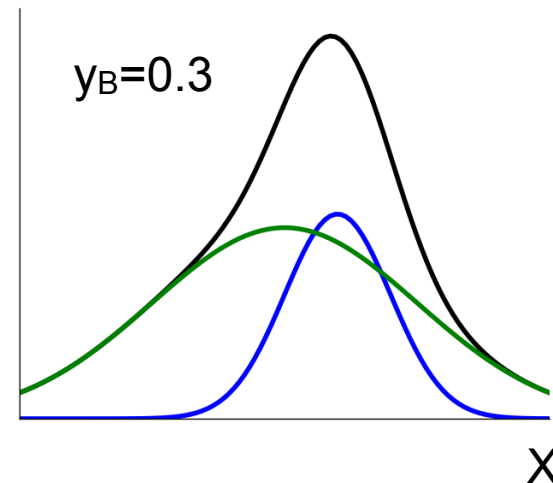
X

$$h_{A,i} = y_A h_{1,i} + (1 - y_A) h_{0,i}$$

$$h_{B,i} = y_B h_{1,i} + (1 - y_B) h_{0,i}$$

- Given two independent unlabeled data samples, and the corresponding proportion of signal, we can extract the signal and background distributions.
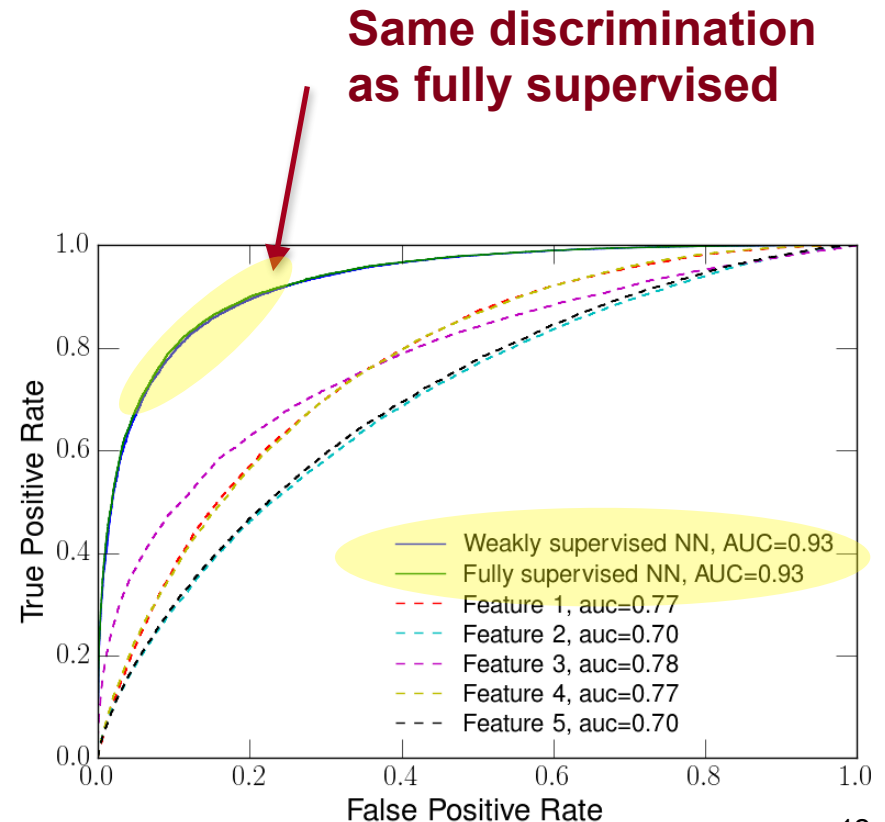—> build Likelihood Ratio discriminant.

11

# Weak supervision

- The analytic approach requires binning and becomes quickly unmanageable as the feature space grows.

- ML approach directly looks for discriminant, without extracting explicitly n-dimensional feature distributions for S and B.
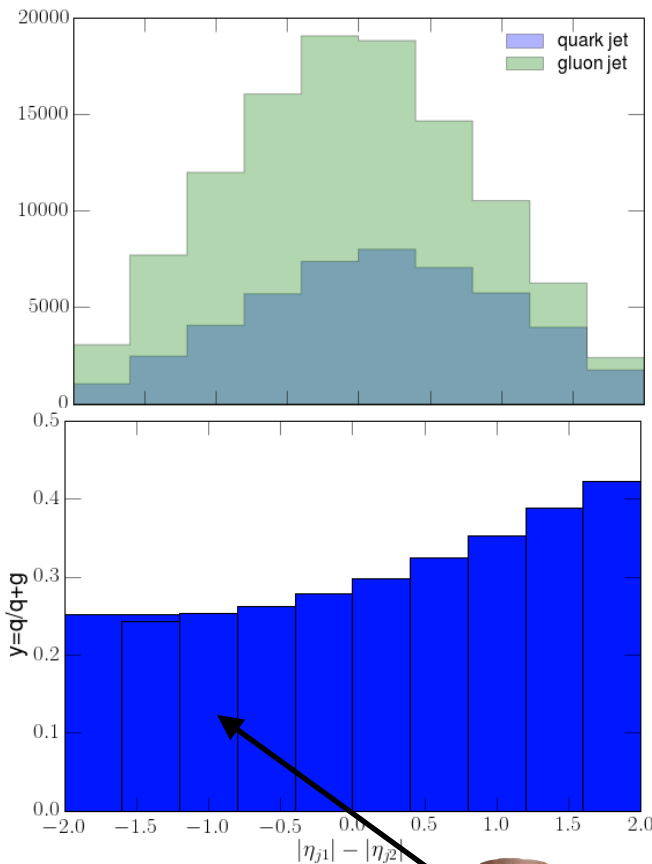
$$f_{\text{full}} = \text{argmin}_{f':\mathbb{R}^n \to \{0,1\}} \sum_{i=1}^{N} \ell \left( f'(x_i) - t_i \right)$$

$$f_{\text{weak}} = \text{argmin}_{f':\mathbb{R}^n \to [0,1]} \ell \left( \sum_{i=1}^{N} \frac{f'(x_i)}{N} - y \right)$$
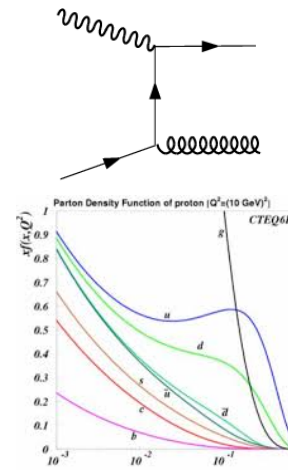
**Same discrimination as fully supervised**



ROC curve legend:
- Weakly supervised NN, AUC=0.93
- Fully supervised NN, AUC=0.93
- Feature 1, auc=0.77
- Feature 2, auc=0.70
- Feature 3, auc=0.78
- Feature 4, auc=0.77
- Feature 5, auc=0.70

x-axis: False Positive Rate
y-axis: True Positive Rate

12

# Weak supervision - q/g tagging

**$|\eta_{j1}| - |\eta_{j2}|$ in dijet events**



quark jet
gluon jet
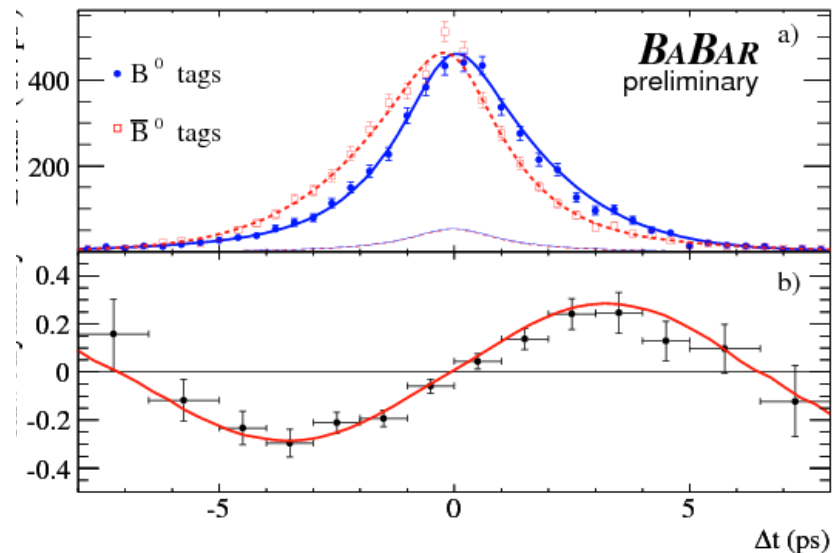
Leverage precise description of ME and PDF (MC/theory) to extract discrimination from soft QCD features (from data!)



1/4 quarks

Each bin is a "barrel" of jets with known proportion



Fully supervised NN, AUC=0.79
Weakly supervised NN, AUC=0.79
n, AUC=0.76
w, AUC=0.78
f0, AUC=0.77

13

# Summary

- **Weak supervision** is a new paradigm leveraging the **class proportions** in high-level observables in order to use **unlabeled data** to extract **discriminating information** from poorly modeled or unknown **low-level observables**.

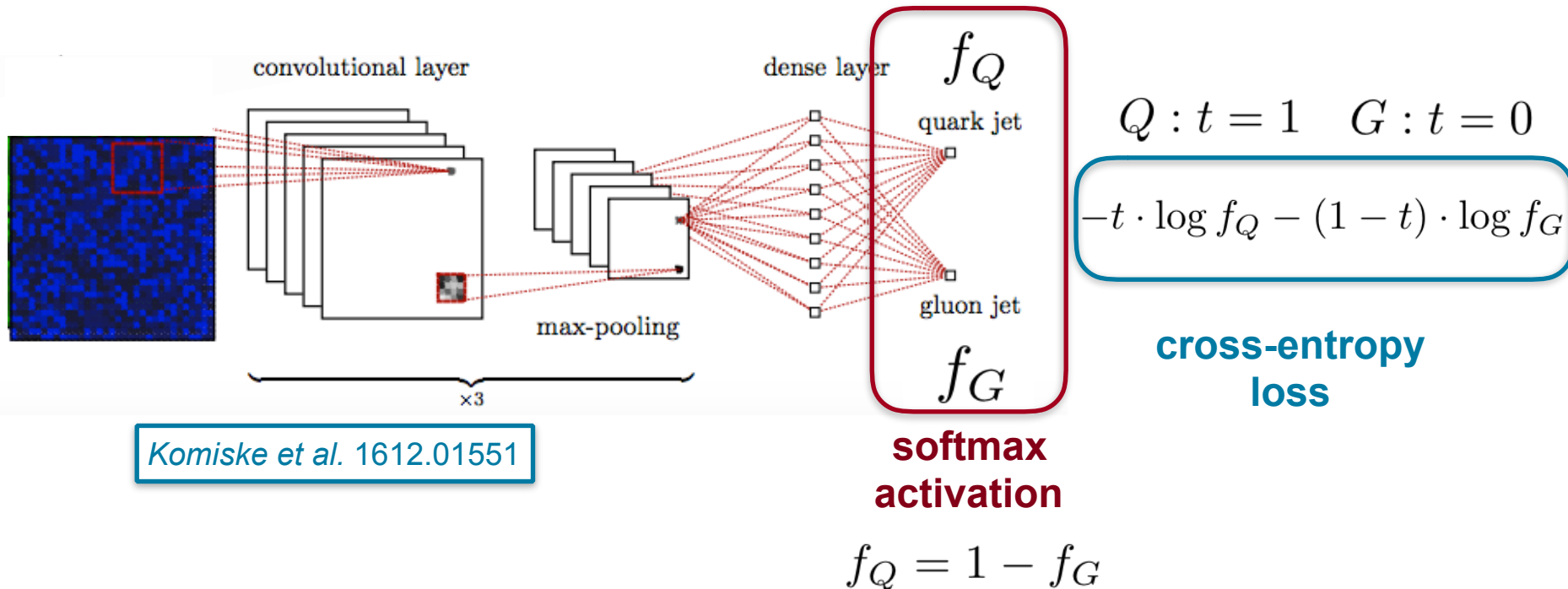- Multiple potential applications in HEP



ATLAS-CONF-2016-055/

SLAC-PUB-13402

14

# Next step: scaling to higher dimensionality

**SLAC**

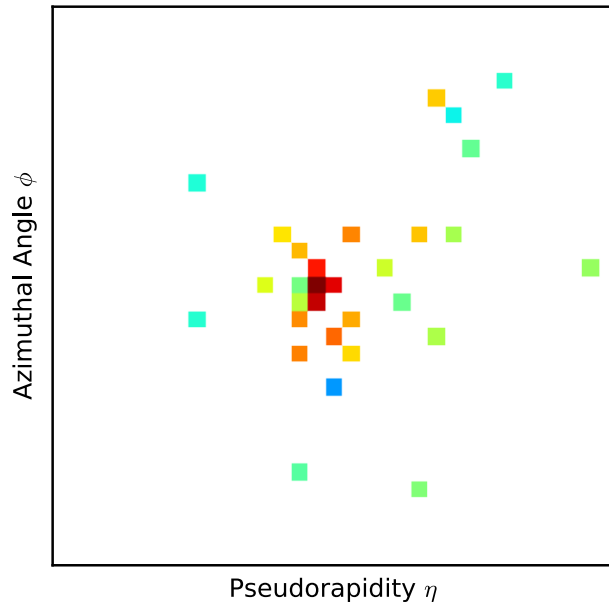Quark/gluon jet tagging with jet images (grayscale) and CNN

Fully supervised network:



*Komiske et al.* 1612.01551

$$Q : t = 1 \quad G : t = 0$$

$$-t \cdot \log f_Q - (1 - t) \cdot \log f_G$$

**cross-entropy loss**

**softmax activation**
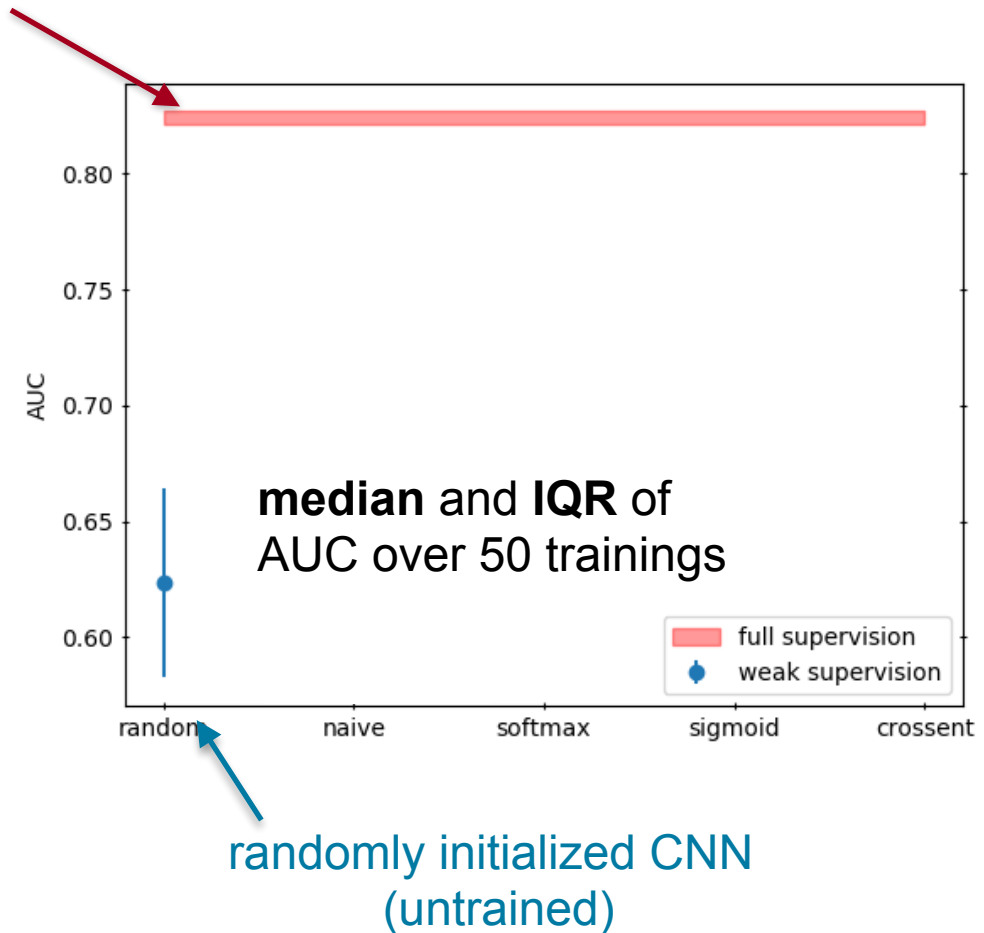
$$f_Q = 1 - f_G$$

First look at weak supervision on same architecture in "ideal" conditions:
**50 samples** with proportions in **[0,1]** (regularly spaced)
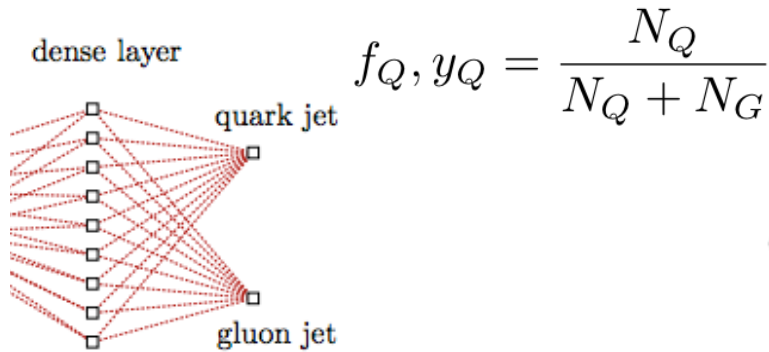
# Jet image + weak supervision

fully supervised CNN



Azimuthal Angle $\phi$

Pseudorapidity $\eta$

33x33=1089 input features

**median** and **IQR** of
AUC over 50 trainings

full supervision
weak supervision

random    naive    softmax    sigmoid    crossent

randomly initialized CNN
(untrained)

# Jet image + weak supervision
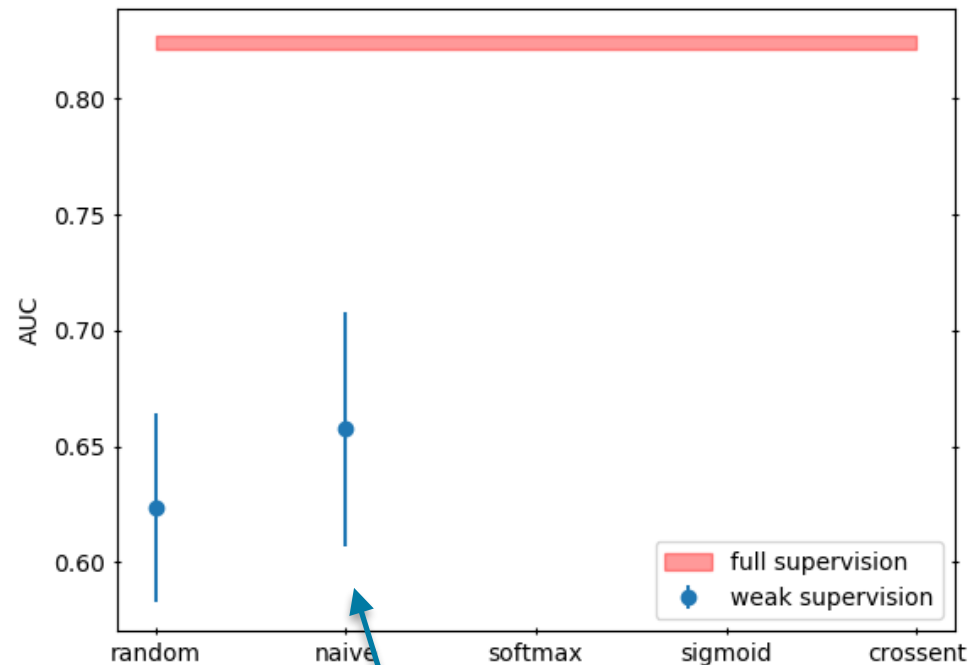
SLAC

dense layer

quark jet

gluon jet

$$f_Q, y_Q = \frac{N_Q}{N_Q + N_G}$$
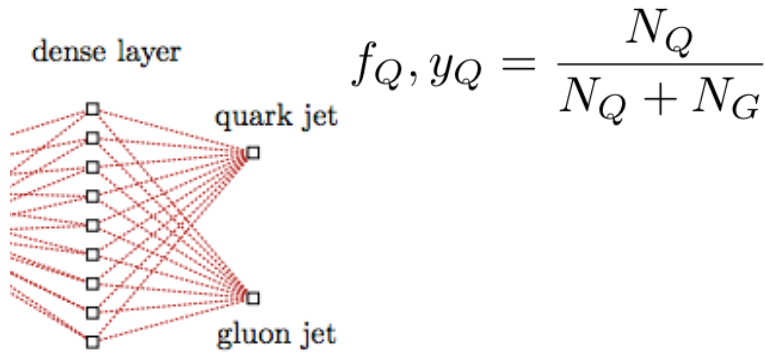
$$f_G, y_G = 1 - y_Q$$

$$L = \left( \sum \frac{f_Q}{N} - y_Q \right)^2$$

less constraint for "gluon" weights
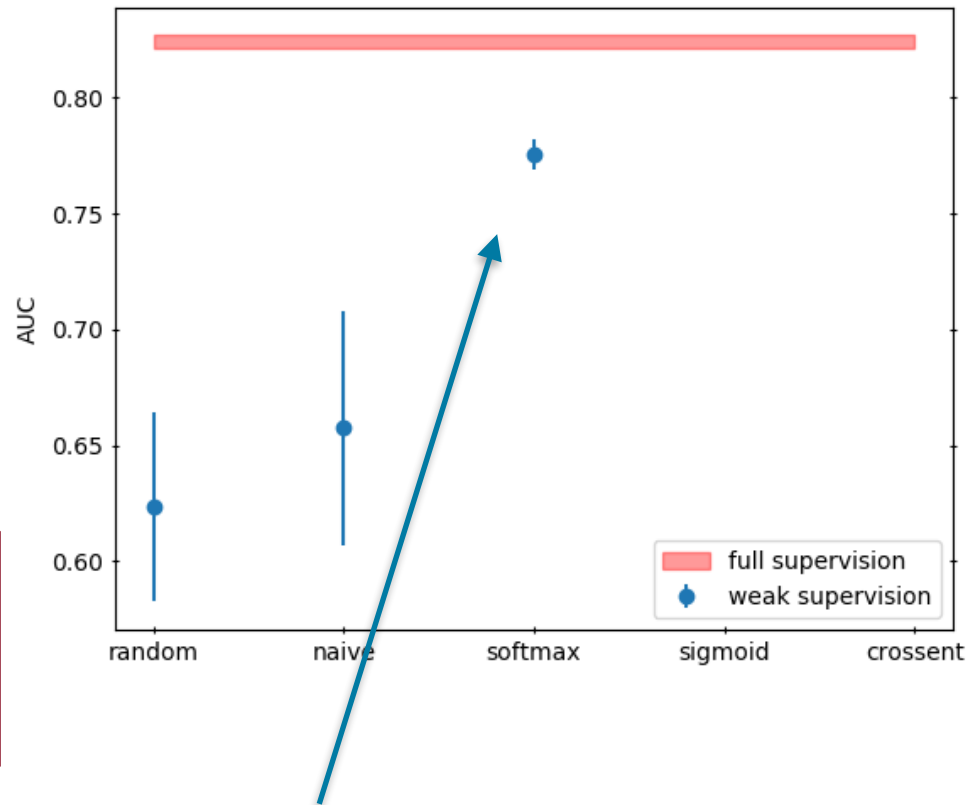(asymmetric gradient)



naive squared loss

# Jet image + weak supervision

SLAC



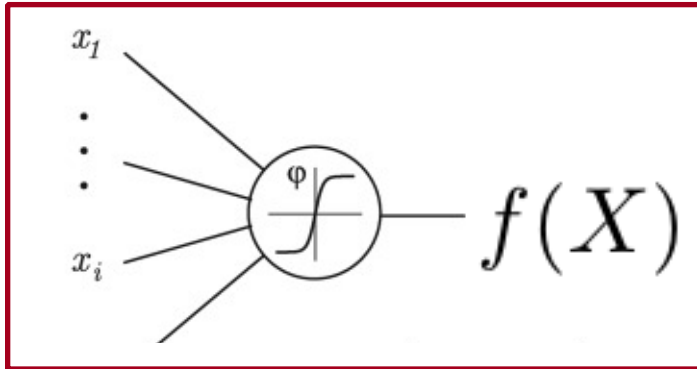$$f_Q, y_Q = \frac{N_Q}{N_Q + N_G}$$

$$f_G, y_G = 1 - y_Q$$

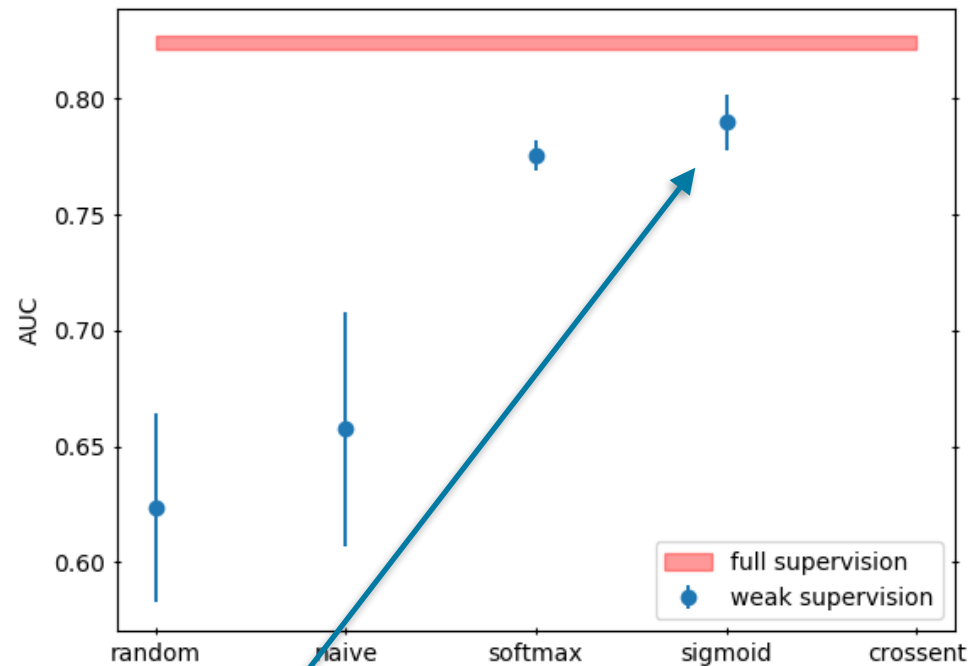$$L = \left( \sum \frac{f_Q}{N} - y_Q \right)^2 + \left( \sum \frac{f_G}{N} - y_G \right)^2$$

symmetric squared loss with softmax activation
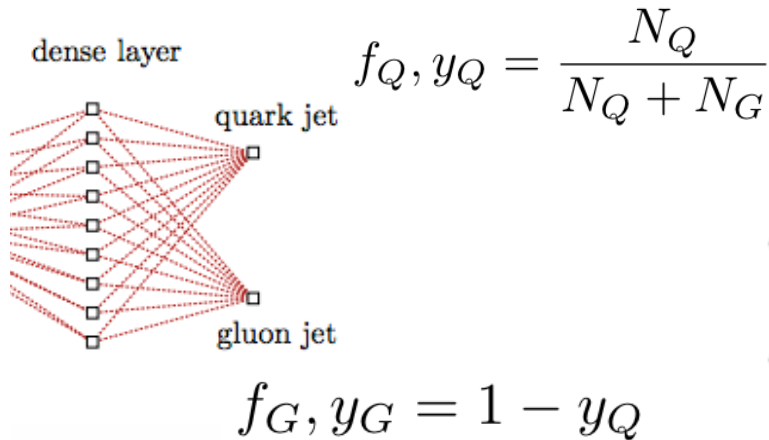
18

# Jet image + weak supervision

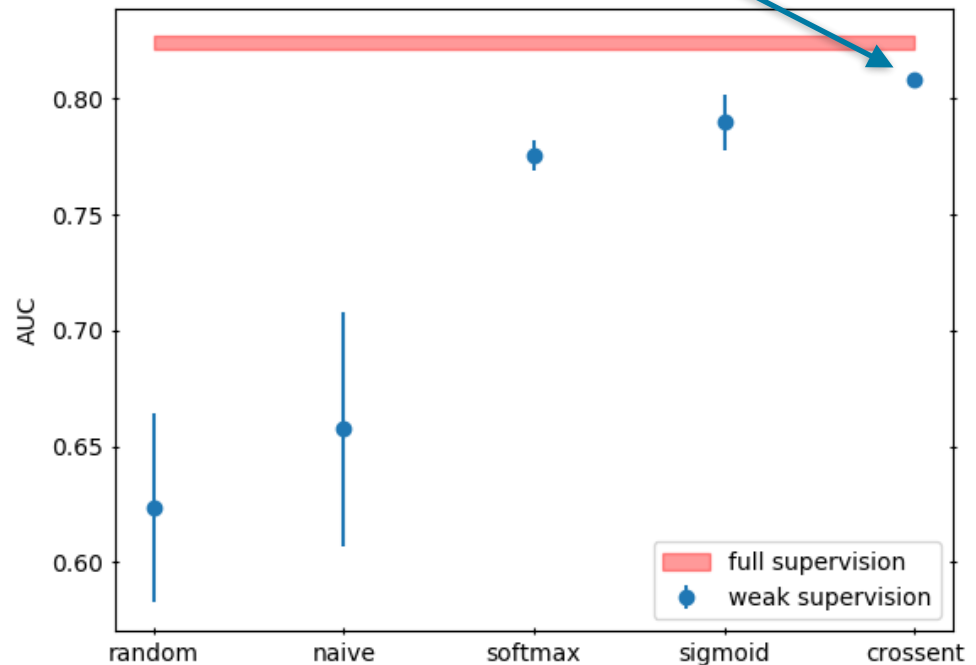SLAC



$$L = \left( \sum \frac{f}{N} - y_Q \right)^2$$

squared loss with sigmoid activation

19

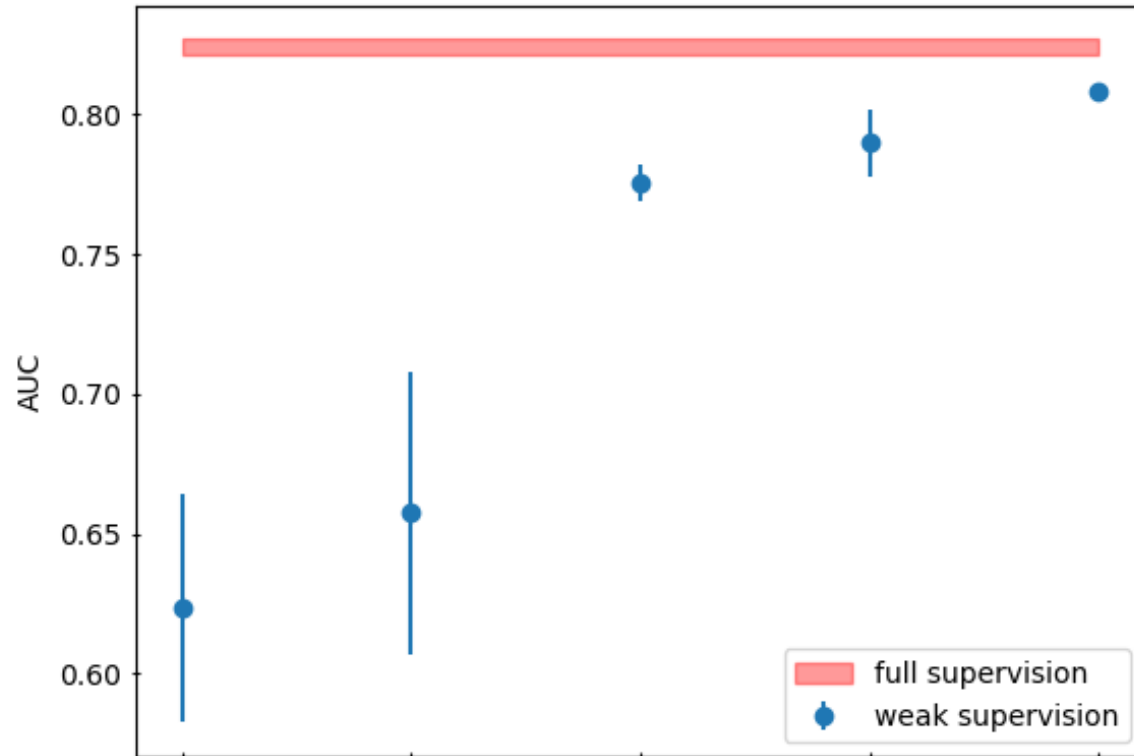# Jet image + weak supervision

"weak cross-entropy" with softmax activation

$$f_Q, y_Q = \frac{N_Q}{N_Q + N_G}$$

$$f_G, y_G = 1 - y_Q$$



$$L = -y_Q \log\left(\sum \frac{f_Q}{N} / y_Q\right) - y_G \log\left(\sum \frac{f_G}{N} / y_G\right)$$

20

# Jet image + weak supervision

SLAC



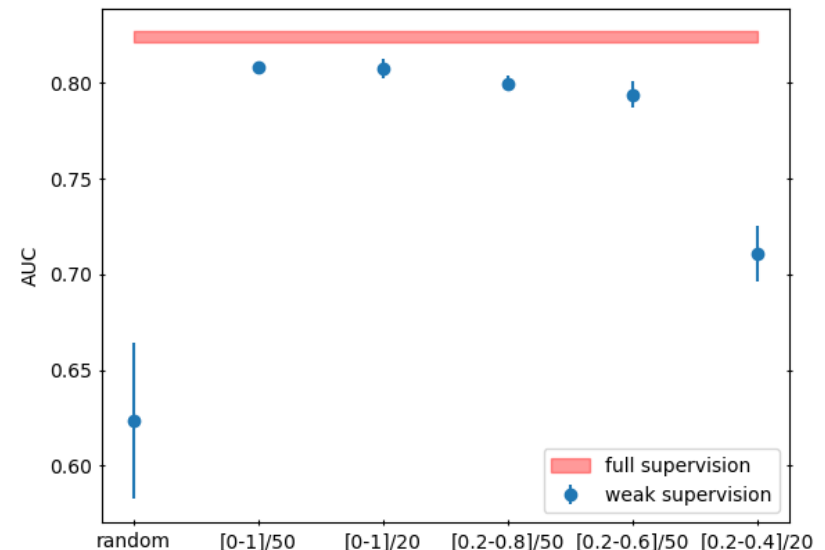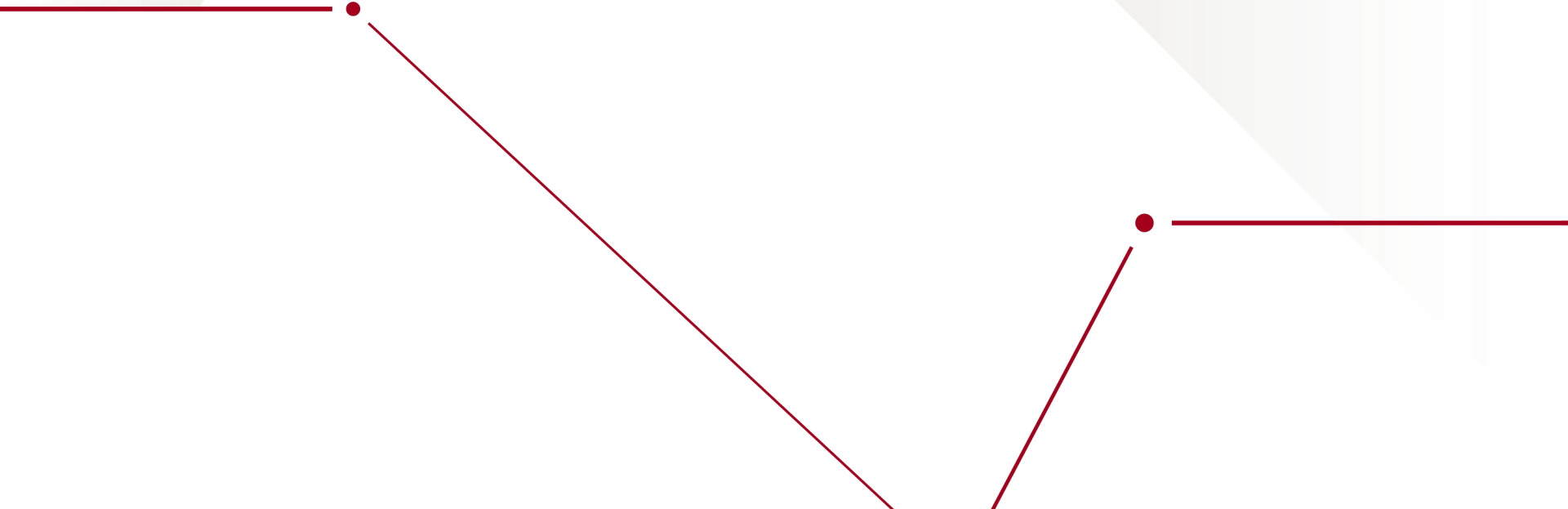| Loss | N/A | square | square symmetric | square symmetric | cross-entropy |
|---|---|---|---|---|---|
| Activation | N/A | softmax | softmax | sigmoid | softmax |

# Conclusions and next steps

- First implementation of **weak supervision+CNN** shows promising results for jet image classification with **unlabeled training data**.

- Careful choice for activation and loss function provide important handles to close gap wrt full supervision performance.

- Plan to investigate impact of size and structure of training data

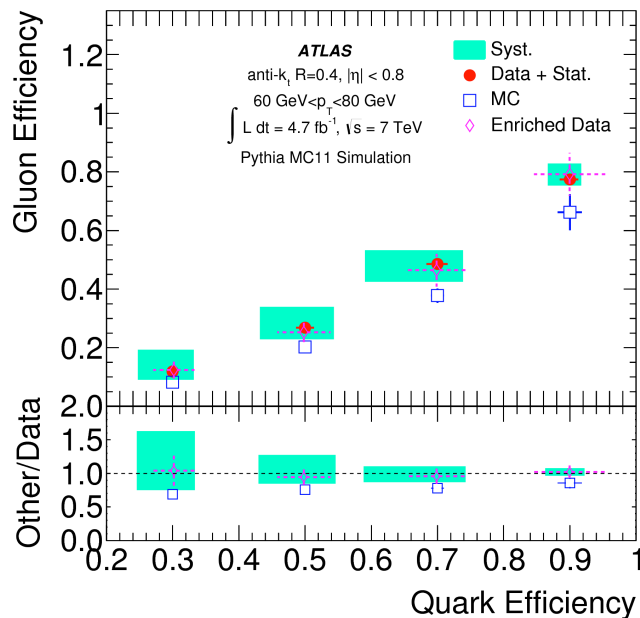- Architecture choices possibly play a role (e.g. "wider" networks)

# References

- Jet-Images: Computer Vision Inspired Techniques for Jet Tagging - https://arxiv.org/abs/1407.5675
- Jet-Images — Deep Learning Edition - https://arxiv.org/abs/1511.05190
- Light-quark and gluon jet discrimination in pp collisions at $\sqrt{s}$=7 TeV with the ATLAS detector - https://arxiv.org/abs/1405.6583
- Weakly Supervised Classification in High Energy Physics - https://arxiv.org/abs/1702.00414
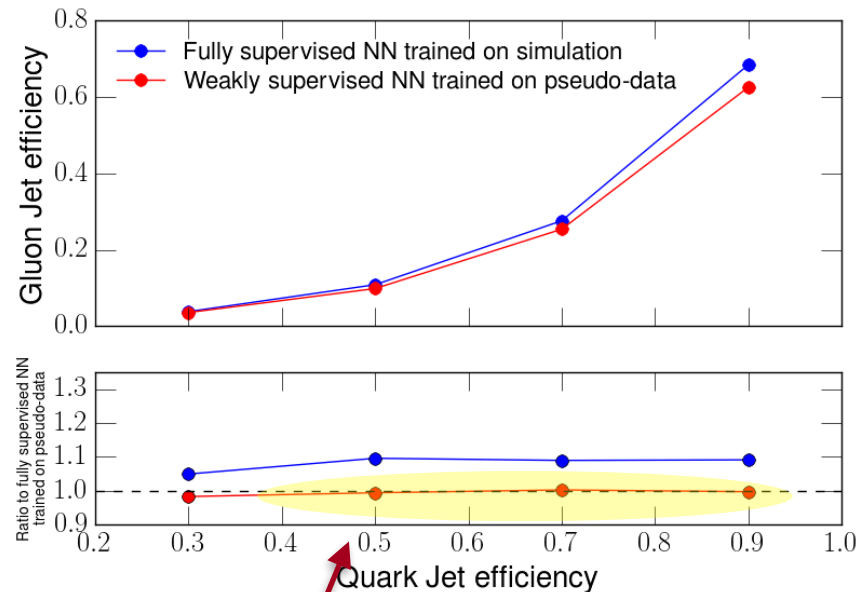
Backup

# Weak supervision

- Weak supervision allows training directly on data
- Learns only <u>real</u> features, from being exposed to discriminant features in data.
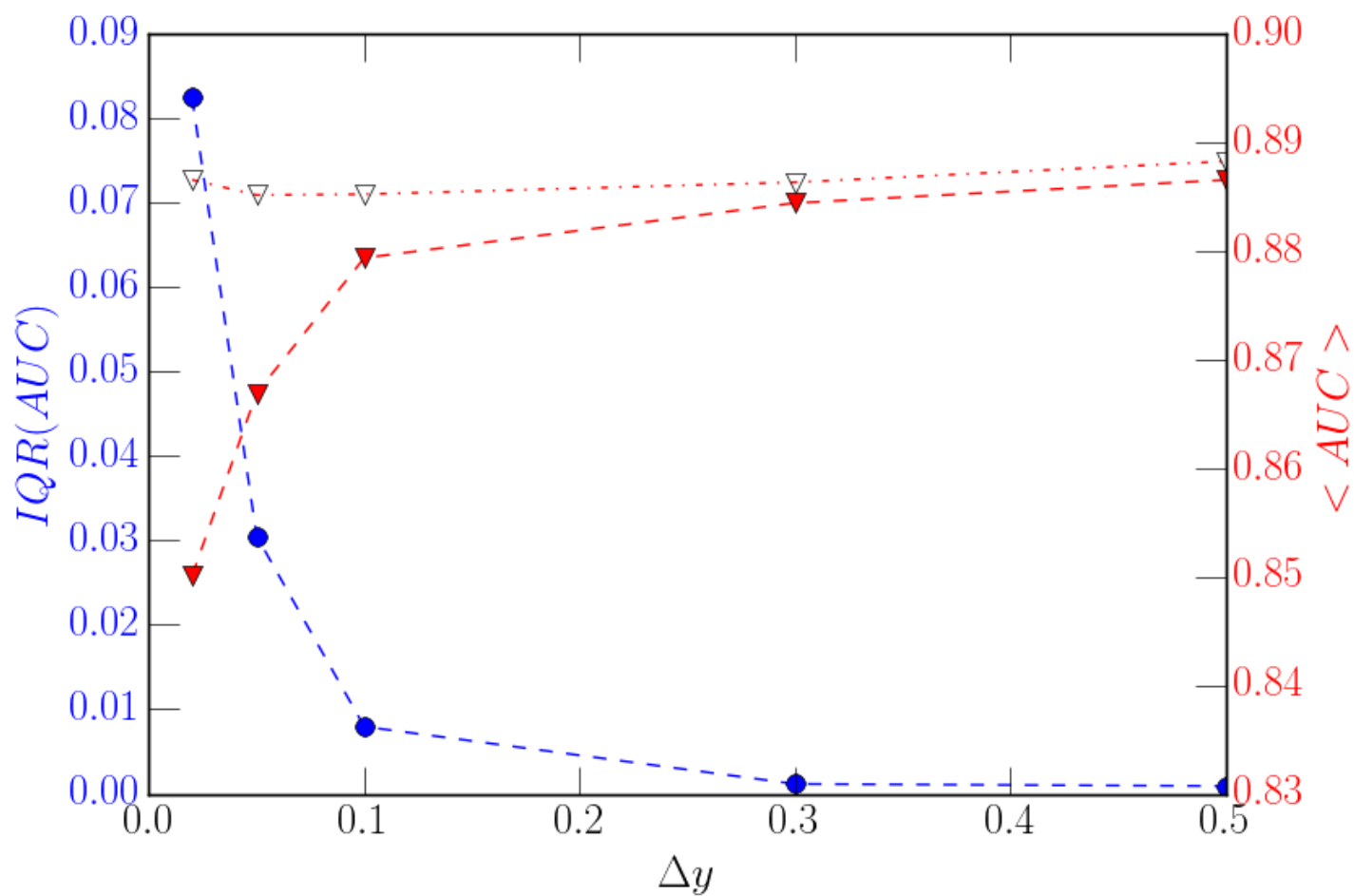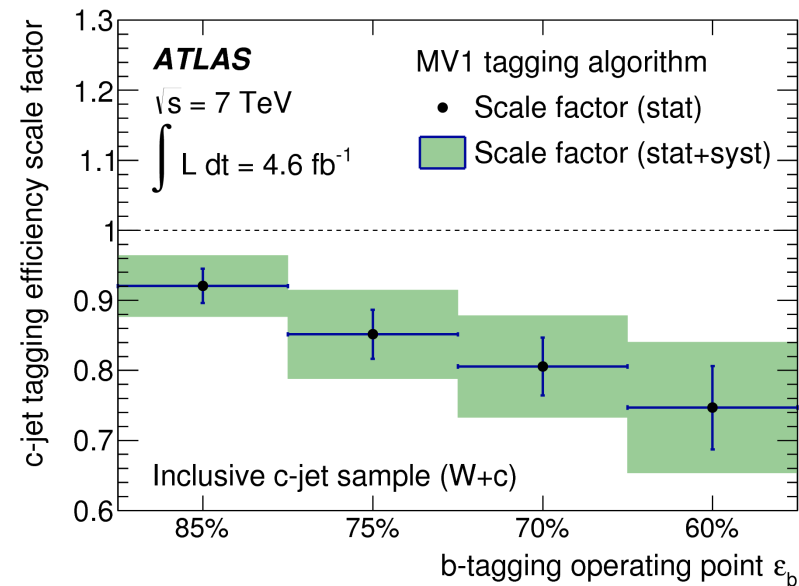


1405.6583

**Same performance as ideal classifier, trained on labeled data**

# Stability

# data-MC SFs

2016 JINST 11 P04008

Cumulative data-simulation scale factor - CMS Tagger, CMS Combined Tagger

| $|\eta| < 1.0$ | | | |
|---|---|---|---|
| Selection | MADGRAPH | POWHEG | MC@NLO |
| CMS Tagger WP0 | $0.985 \pm 0.073$ | $1.173 \pm 0.092$ | $1.033 \pm 0.081$ |
| CMS Combined Tagger WP3 | $0.891 \pm 0.118$ | $1.063 \pm 0.146$ | $0.933 \pm 0.129$ |

| $1.0 < |\eta| < 2.4$ | | | |
|---|---|---|---|
| Selection | MADGRAPH | POWHEG | MC@NLO |
| CMS Tagger WP0 | $0.644 \pm 0.100$ | $0.704 \pm 0.110$ | $0.768 \pm 0.118$ |
| CMS Combined Tagger WP3 | $0.685 \pm 0.199$ | $0.906 \pm 0.277$ | $0.802 \pm 0.230$ |

CMS-PAS-JME-13-007

# Jet image + weak supervision